

ADDENDUM ON THE SCORING OF GAUSSIAN DIRECTED ACYCLIC GRAPHICAL MODELS

BY JACK KUIPERS^{*}, GIUSI MOFFA^{*} AND DAVID HECKERMAN[†]

Regensburg University^{} and Microsoft Research[†]*

We provide a correction to the expression for scoring Gaussian directed acyclic graphical models derived in Geiger and Heckerman (2002) [*Annals of Statistics* **20** 1414-1440] and discuss how to evaluate the score efficiently.

Gaussian directed acyclic graph (DAG) models represent a particular type of Bayesian networks where the node variables are assumed to come from a multivariate Gaussian distribution. The Bayesian Gaussian equivalent (BGe) score was introduced in Geiger and Heckerman (1994); Heckerman and Geiger (1995); Geiger and Heckerman (2002) for learning such networks.

We follow the same notation as Geiger and Heckerman (2002) in considering DAG models m with n nodes corresponding to the set of variables $\mathbf{X} = \{X_1, \dots, X_n\}$. Let m^h be the model hypothesis that the true distribution of \mathbf{X} is faithful to the DAG model m , meaning that it satisfies only and all the conditional independencies encoded by the DAG. For a complete random data sample $d = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with N observations and a complete DAG model m_c , the marginal likelihood is (Geiger and Heckerman, 2002, Theorem 2)

$$(1) \quad p(d \mid m^h) = \prod_{i=1}^n \frac{p(d^{\mathbf{Pa}_i \cup \{X_i\}} \mid m_c^h)}{p(d^{\mathbf{Pa}_i} \mid m_c^h)}$$

where \mathbf{Pa}_i are the parent variables of the vertex i and $d^{\mathbf{Y}}$ is the data restricted to the coordinates in $\mathbf{Y} \subseteq \mathbf{X}$. Different DAGs which encode the same set of conditional independence are said to belong to an equivalence class. Along with ensuring that all DAGs in the same equivalence class are scored equally, the modularity of the score allows the steps in structure MCMC (Madigan and York, 1995) to be evaluated much more efficiently. Order MCMC (Friedman and Koller, 2003, on the related space of triangular matrices) as well as the edge reversal move of Grzegorzcyk and Husmeier (2008) would not be possible without it.

MSC 2010 subject classifications: 62-07, 62F15, 62H99

Keywords and phrases: Gaussian DAG models, Bayesian network learning, BGe score

For Gaussian DAG models, the likelihood is a multivariate normal distribution with mean $\boldsymbol{\mu}$ and precision matrix W . The need for global parameter independence, so that the expression of the score in (1) holds, implies that the prior distribution of $(\boldsymbol{\mu}, W)$ must be normal-Wishart (Geiger and Heckerman, 2002). The parameter $\boldsymbol{\mu}$ is taken to be normally distributed with mean $\boldsymbol{\nu}$ and precision matrix $\alpha_\mu W$, for $\alpha_\mu > 0$. W is Wishart distributed with positive definite parametric matrix T (the inverse of the scale matrix) and degrees of freedom α_w , with $\alpha_w > n - 1$. As detailed in the Supplement, one finds

$$(2) \quad p(d^{\mathbf{Y}} \mid m_c^h) = \left(\frac{\alpha_\mu}{N + \alpha_\mu} \right)^{\frac{l}{2}} \frac{\Gamma_l \left(\frac{N + \alpha_w - n + l}{2} \right)}{\pi^{\frac{lN}{2}} \Gamma_l \left(\frac{\alpha_w - n + l}{2} \right)} \frac{|T_{\mathbf{Y}\mathbf{Y}}|^{\frac{\alpha_w - n + l}{2}}}{|R_{\mathbf{Y}\mathbf{Y}}|^{\frac{N + \alpha_w - n + l}{2}}}$$

where l is the size of \mathbf{Y} , $A_{\mathbf{Y}\mathbf{Y}}$ means selecting the rows and columns corresponding to \mathbf{Y} of a matrix A ,

$$(3) \quad \Gamma_l \left(\frac{x}{2} \right) = \pi^{\frac{n(n-1)}{4}} \prod_{j=1}^l \Gamma \left(\frac{x + 1 - j}{2} \right)$$

is the multivariate Gamma function and

$$(4) \quad R = T + S_N + \frac{N\alpha_w}{(N + \alpha_w)} (\boldsymbol{\nu} - \bar{\mathbf{x}}) (\boldsymbol{\nu} - \bar{\mathbf{x}})^T$$

is the posterior parametric matrix involving

$$(5) \quad \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad S_N = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

the sample mean and sample variance multiplied by $(N - 1)$.

The result in (2) is identical to equation (18) of Geiger and Heckerman (2002), once some factors are cancelled, apart from the manner in which the matrix elements are chosen. The result in Geiger and Heckerman (2002) replaces the $T_{\mathbf{Y}\mathbf{Y}}$ and $R_{\mathbf{Y}\mathbf{Y}}$ by $T_{\mathbf{Y}}$ and $R_{\mathbf{Y}}$, where $A_{\mathbf{Y}} = ((A^{-1})_{\mathbf{Y}\mathbf{Y}})^{-1}$. Inverting the matrices before the elements are selected and then inverting again (as in Geiger and Heckerman, 2002) we found non-consistent behaviour on simulated data.

We may further compare to equation (24) of Heckerman and Geiger (1995), which with the current notation becomes

$$(6) \quad p(d^{\mathbf{Y}} \mid m_c^h) = \left(\frac{\alpha_\mu}{N + \alpha_\mu} \right)^{\frac{l}{2}} \frac{\Gamma_l \left(\frac{N + \alpha_w}{2} \right)}{\pi^{\frac{lN}{2}} \Gamma_l \left(\frac{\alpha_w}{2} \right)} \frac{|T_{\mathbf{Y}\mathbf{Y}}|^{\frac{\alpha_w}{2}}}{|R_{\mathbf{Y}\mathbf{Y}}|^{\frac{N + \alpha_w}{2}}}$$

while incorrectly defining the S_N in the R in (4) as the sample variance. However, the same terminology, with the correct formula for S_N , is used in Geiger and Heckerman (1994) whose equation (15) is otherwise identical to (6) aside from having π replaced by 2π .

The difference in the powers of the determinants between (2) and (6) could lead to a subtle, and hard to predict, change in the scores. There is also the same loss of l -dependence in the arguments of the multivariate gamma functions. The ratio of gamma functions for each node now actually decreases with l while the ratio from (2) increases instead. Using (6) instead of (2) effectively penalises each node with l parents by a factor $\sim N^l$, giving a substantial bias towards sparse DAGs. This bias is likely to be present in early works implementing the score of Heckerman and Geiger (1995) and possibly remains in legacy code.

SUPPLEMENTARY MATERIAL

Supplement: Deriving and simplifying the BGe score

We detail the steps used to derive (2) and simplify the ratios appearing in (1) to improve the numerical evaluation of the score.

A. Deriving the score. We first proceed through the steps needed to derive the Bayesian Gaussian equivalent (BGe) score as introduced in Geiger and Heckerman (1994); Heckerman and Geiger (1995); Geiger and Heckerman (2002). We use the same notation as Geiger and Heckerman (2002) and let n be the dimension of the variables \mathbf{X} , N the number of observations of \mathbf{x} in a dataset $d = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. For a Gaussian directed acyclic graphical (DAG) model, we wish to evaluate the terms in the product

$$(A.1) \quad p(d \mid m^h) = \prod_{i=1}^n \frac{p(d^{\mathbf{Pa}_i \cup \{X_i\}} \mid m_c^h)}{p(d^{\mathbf{Pa}_i} \mid m_c^h)}$$

where \mathbf{Pa}_i are the parent variables of the vertex i and $d^{\mathbf{Y}}$ is the data restricted to the coordinates in $\mathbf{Y} \subseteq \mathbf{X}$. In the following we leave out the explicit dependence on the model hypothesis m^h and the complete data model m_c^h in the formulae.

A.1. Definitions. The data is assumed to be normally distributed with mean $\boldsymbol{\mu}$ and precision matrix W

$$(A.2) \quad \begin{aligned} p(d \mid \boldsymbol{\mu}, W) &= \frac{|W|^{\frac{N}{2}}}{(2\pi)^{\frac{nN}{2}}} e^{-\frac{1}{2} \sum_{i=1}^N (\boldsymbol{\mu} - \mathbf{x}_i)^T W (\boldsymbol{\mu} - \mathbf{x}_i)} \\ &= \frac{|W|^{\frac{N}{2}}}{(2\pi)^{\frac{nN}{2}}} e^{-\frac{1}{2} \text{Tr}[(\sum_{i=1}^N (\boldsymbol{\mu} - \mathbf{x}_i)(\boldsymbol{\mu} - \mathbf{x}_i)^T] W)} \end{aligned}$$

where we reordered the exponent using the cyclicity of the trace. The prior on W is taken to be a Wishart distribution, $W \sim \mathcal{W}_n(T^{-1}, \alpha_w)$, where $\alpha_w > n-1$ is the degrees of freedom and T is the positive definite parametric matrix. We follow the standard practice, as in [Press \(1982\)](#), of using the scale matrix (the inverse of the parametric matrix) in the argument of the distribution, which has the following form

$$(A.3) \quad p(W) = \frac{|W|^{\frac{\alpha_w - n - 1}{2}}}{Z_{\mathcal{W}}(n, T, \alpha_w)} e^{-\frac{1}{2} \text{Tr}[TW]}$$

with normalising constant

$$(A.4) \quad Z_{\mathcal{W}}(n, T, \alpha_w) = \frac{2^{\frac{\alpha_w n}{2}} \Gamma_n\left(\frac{\alpha_w}{2}\right)}{|T|^{\frac{\alpha_w}{2}}}$$

involving the multivariate gamma function

$$(A.5) \quad \Gamma_n\left(\frac{\alpha_w}{2}\right) = \pi^{\frac{n(n-1)}{4}} \prod_{j=1}^n \Gamma\left(\frac{\alpha_w + 1 - j}{2}\right)$$

Next assume a normal prior on $\boldsymbol{\mu}$ with mean $\boldsymbol{\nu}$ and precision matrix $\alpha_\mu W$, with $\alpha_\mu > 0$,

$$(A.6) \quad \begin{aligned} p(\boldsymbol{\mu}|W) &= \frac{1}{Z_{\mathcal{N}}(n, W, \alpha_\mu)} e^{-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\nu})^T \alpha_\mu W (\boldsymbol{\mu} - \boldsymbol{\nu})} \\ &= \frac{1}{Z_{\mathcal{N}}(n, W, \alpha_\mu)} e^{-\frac{1}{2} \text{Tr}[(\boldsymbol{\mu} - \boldsymbol{\nu})(\boldsymbol{\mu} - \boldsymbol{\nu})^T \alpha_\mu W]} \end{aligned}$$

again using the cyclicity of the trace and where the normalising constant is

$$(A.7) \quad Z_{\mathcal{N}}(n, W, \alpha_\mu) = \frac{(2\pi)^{\frac{n}{2}}}{(\alpha_\mu)^{\frac{n}{2}} |W|^{\frac{1}{2}}}$$

Jointly $\boldsymbol{\mu}, W$ therefore follow a normal-Wishart prior distribution

$$(A.8) \quad p(\boldsymbol{\mu}, W) = \frac{|W|^{\frac{\alpha_w - n}{2}}}{Z(n, \alpha_\mu, T, \alpha_w)} e^{-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\nu})^T \alpha_\mu W (\boldsymbol{\mu} - \boldsymbol{\nu})} e^{-\frac{1}{2} \text{Tr}[TW]}$$

with

$$(A.9) \quad Z(n, \alpha_\mu, T, \alpha_w) = \frac{2^{\frac{(\alpha_w + 1)n}{2}} \pi^{\frac{n}{2}} \Gamma_n\left(\frac{\alpha_w}{2}\right)}{(\alpha_\mu)^{\frac{n}{2}} |T|^{\frac{\alpha_w}{2}}}$$

A.2. *The posterior distribution.* The first step is to calculate the posterior probability $p(\boldsymbol{\mu}, W|d) = p(d, \boldsymbol{\mu}, W)/P(d)$ of the parameters $\boldsymbol{\mu}, W$ given the data d . Start by writing the joint probability in the form $p(d, \boldsymbol{\mu}, W) = p(d|\boldsymbol{\mu}, W)p(\boldsymbol{\mu}, W)$ and taking the log

$$\begin{aligned} \ln p(d, \boldsymbol{\mu}, W) = & \dots \\ & -\frac{1}{2} \text{Tr} \left[\left(T + \alpha_\mu (\boldsymbol{\mu} - \boldsymbol{\nu}) (\boldsymbol{\mu} - \boldsymbol{\nu})^T + \sum_{i=1}^N (\boldsymbol{\mu} - \mathbf{x}_i) (\boldsymbol{\mu} - \mathbf{x}_i)^T \right) W \right] \end{aligned} \quad (\text{A.10})$$

Define

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad S_N = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (\text{A.11})$$

so that S_N is the sample variance multiplied by $(N-1)$. By expanding and comparing terms, one has

$$\sum_{i=1}^N (\boldsymbol{\mu} - \mathbf{x}_i) (\boldsymbol{\mu} - \mathbf{x}_i)^T = S_N + N (\boldsymbol{\mu} - \bar{\mathbf{x}}) (\boldsymbol{\mu} - \bar{\mathbf{x}})^T \quad (\text{A.12})$$

and also

$$\begin{aligned} & N (\boldsymbol{\mu} - \bar{\mathbf{x}}) (\boldsymbol{\mu} - \bar{\mathbf{x}})^T + \alpha_\mu (\boldsymbol{\mu} - \boldsymbol{\nu}) (\boldsymbol{\mu} - \boldsymbol{\nu})^T \\ &= (N + \alpha_\mu) (\boldsymbol{\mu} - \boldsymbol{\nu}') (\boldsymbol{\mu} - \boldsymbol{\nu}')^T + \frac{N\alpha_\mu}{(N + \alpha_\mu)} (\bar{\mathbf{x}} - \boldsymbol{\nu}) (\bar{\mathbf{x}} - \boldsymbol{\nu})^T \end{aligned} \quad (\text{A.13})$$

with

$$\boldsymbol{\nu}' = \frac{N\bar{\mathbf{x}} + \alpha_\mu \boldsymbol{\nu}}{(N + \alpha_\mu)} \quad (\text{A.14})$$

Further define

$$R = T + S_N + \frac{N\alpha_\mu}{(N + \alpha_\mu)} (\bar{\mathbf{x}} - \boldsymbol{\nu}) (\bar{\mathbf{x}} - \boldsymbol{\nu})^T \quad (\text{A.15})$$

then we may rewrite the joint probability from (A.10) as

$$\ln p(d, \boldsymbol{\mu}, W) = \dots - \frac{1}{2} \text{Tr} \left[\left(R + (N + \alpha_\mu) (\boldsymbol{\mu} - \boldsymbol{\nu}') (\boldsymbol{\mu} - \boldsymbol{\nu}')^T \right) W \right] \quad (\text{A.16})$$

or explicitly

$$\begin{aligned} & (\text{A.17}) \\ p(d, \boldsymbol{\mu}, W) = & \frac{|W|^{\frac{N+\alpha_w-n}{2}}}{(2\pi)^{\frac{nN}{2}} Z(n, \alpha_\mu, T, \alpha_w)} e^{-\frac{1}{2}(\boldsymbol{\mu}-\boldsymbol{\nu}')^T (N+\alpha_\mu) W (\boldsymbol{\mu}-\boldsymbol{\nu}')} e^{-\frac{1}{2} \text{Tr}[RW]} \end{aligned}$$

In order to obtain the marginal $P(d)$ the above expression needs to be integrated with respect to $\boldsymbol{\mu}, W$. Comparing to (A.8) we can see that the functional form inside the integral corresponds to that of a normal-Wishart distribution with the following transformations

$$\begin{aligned} \alpha_\mu &\rightarrow N + \alpha_\mu \\ \alpha_w &\rightarrow N + \alpha_w \\ \boldsymbol{\nu} &\rightarrow \boldsymbol{\nu}' \\ T &\rightarrow R \end{aligned} \tag{A.18}$$

Therefore its integral is given by the normalization constant, leading to (A.19)

$$p(d) = \frac{Z(n, N + \alpha_\mu, R, N + \alpha_w)}{(2\pi)^{\frac{nN}{2}} Z(n, \alpha_\mu, T, \alpha_w)} = \left(\frac{\alpha_\mu}{N + \alpha_\mu} \right)^{\frac{n}{2}} \frac{\Gamma_n\left(\frac{N + \alpha_w}{2}\right)}{\pi^{\frac{nN}{2}} \Gamma_n\left(\frac{\alpha_w}{2}\right)} \frac{|T|^{\frac{\alpha_\mu}{2}}}{|R|^{\frac{N + \alpha_\mu}{2}}}$$

This formula agrees with both equation (24) of Heckerman and Geiger (1995) and equation (18) of Geiger and Heckerman (2002) in the case when the full data is considered rather than a node specific subset.

From the above it also follows that the posterior distribution of the parameters $\boldsymbol{\mu}, W$ remains normal-Wishart

$$p(\boldsymbol{\mu}, W|d) = \frac{|W|^{\frac{N + \alpha_w - n}{2}}}{Z(n, N + \alpha_\mu, R, N + \alpha_w)} e^{-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\nu}')^T (N + \alpha_\mu) W (\boldsymbol{\mu} - \boldsymbol{\nu}')} e^{-\frac{1}{2} \text{Tr}[RW]} \tag{A.20}$$

with parameters as in (A.18).

A.3. Subsets. Finally we need to perform the same steps but when we restrict to a subset \mathbf{Y} of size l of the n coordinates of the data, $\mathbf{Y} \subseteq \mathbf{X}$. These subsets are the data corresponding to either the node being scored and its parents, or its parents alone, and the resulting terms are needed for the complete evaluation of the BGe score in (A.1). We form a mean vector $\boldsymbol{\mu}_{\mathbf{Y}}$ from the components of $\boldsymbol{\mu}$ that are in \mathbf{Y} and similarly partition the matrix W . Denoting the complement of \mathbf{Y} by $\tilde{\mathbf{Y}}$, and with the elements appropriately reordered

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_{\mathbf{Y}} \\ \boldsymbol{\mu}_{\tilde{\mathbf{Y}}} \end{pmatrix}, \quad W = \begin{pmatrix} W_{\mathbf{Y}\mathbf{Y}} & W_{\mathbf{Y}\tilde{\mathbf{Y}}} \\ W_{\tilde{\mathbf{Y}}\mathbf{Y}} & W_{\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}} \end{pmatrix}, \quad W_{\mathbf{Y}\tilde{\mathbf{Y}}} = W_{\tilde{\mathbf{Y}}\mathbf{Y}}^T \tag{A.21}$$

Since the data d is assumed to be normally distributed with mean $\boldsymbol{\mu}$ and precision matrix W , the subset of the data on the \mathbf{Y} coordinates is then

normally distributed with mean $\boldsymbol{\mu}_{\mathbf{Y}}$ and precision matrix $W_{\mathbf{Y}}$ where

$$(A.22) \quad W_{\mathbf{Y}} = \left((W^{-1})_{\mathbf{Y}\mathbf{Y}} \right)^{-1} = W_{\mathbf{Y}\mathbf{Y}} - W_{\mathbf{Y}\tilde{\mathbf{Y}}} W_{\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}}^{-1} W_{\tilde{\mathbf{Y}}\mathbf{Y}}$$

Since for a multivariate normal we just need to keep the corresponding elements of the covariance matrix, we simply invert W , select the elements and then invert again as in (A.22).

Next we need the distribution of $\boldsymbol{\mu}_{\mathbf{Y}}$ given that $\boldsymbol{\mu}$ is normally distributed with mean $\boldsymbol{\nu} = \begin{pmatrix} \boldsymbol{\nu}_{\mathbf{Y}} \\ \boldsymbol{\nu}_{\tilde{\mathbf{Y}}} \end{pmatrix}$ and precision matrix $\alpha_{\mu} W$. Again transforming to the covariance matrix, selecting elements and transforming back leads to a normal distribution with mean $\boldsymbol{\nu}_{\mathbf{Y}}$ and precision matrix $\alpha_{\mu} W_{\mathbf{Y}}$. All we really need then is the distribution of $W_{\mathbf{Y}}$ given that W follows a Wishart distribution with parametric matrix T and degrees of freedom α_w . For this we use Theorem 5.1.4 of Press (1982)

- If $W \sim \mathcal{W}_n(T^{-1}, \alpha_w)$ is Wishart distributed, then

$$(A.23) \quad W_{\mathbf{Y}} \sim \mathcal{W}_l \left((T_{\mathbf{Y}\mathbf{Y}})^{-1}, \alpha_w - n + l \right)$$

where the degrees of freedom has been reduced and we simply select the relevant entries from the parametric matrix.

This result is also included in Theorem 5 of Geiger and Heckerman (2002).

The prior distribution for the subset of the parameters is then

$$(A.24) \quad p(\boldsymbol{\mu}_{\mathbf{Y}}, W_{\mathbf{Y}}) = \frac{|W_{\mathbf{Y}}|^{\frac{\alpha_w - n}{2}}}{Z(l, \alpha_{\mu}, T_{\mathbf{Y}\mathbf{Y}}, \alpha_w - n + l)} \cdot e^{-\frac{1}{2}(\boldsymbol{\mu}_{\mathbf{Y}} - \boldsymbol{\nu}_{\mathbf{Y}})^T \alpha_{\mu} W_{\mathbf{Y}} (\boldsymbol{\mu}_{\mathbf{Y}} - \boldsymbol{\nu}_{\mathbf{Y}})} e^{-\frac{1}{2} \text{Tr}[T_{\mathbf{Y}\mathbf{Y}} W_{\mathbf{Y}}]}$$

while the likelihood for the corresponding subset of data $d^{\mathbf{Y}}$ is

$$(A.25) \quad p(d^{\mathbf{Y}} | \boldsymbol{\mu}_{\mathbf{Y}}, W_{\mathbf{Y}}) = \frac{|W_{\mathbf{Y}}|^{\frac{N}{2}}}{(2\pi)^{\frac{IN}{2}}} e^{-\frac{1}{2} \text{Tr}[(S_N)_{\mathbf{Y}\mathbf{Y}} + N(\boldsymbol{\mu}_{\mathbf{Y}} - \bar{\mathbf{x}}_{\mathbf{Y}})(\boldsymbol{\mu}_{\mathbf{Y}} - \bar{\mathbf{x}}_{\mathbf{Y}})^T] W_{\mathbf{Y}}]}$$

We may follow the same steps as in section A.2 to obtain the joint probability

$$(A.26) \quad p(d^{\mathbf{Y}}, \boldsymbol{\mu}_{\mathbf{Y}}, W_{\mathbf{Y}}) = \frac{|W_{\mathbf{Y}}|^{\frac{N + \alpha_w - n}{2}}}{(2\pi)^{\frac{IN}{2}} Z(l, \alpha_{\mu}, T_{\mathbf{Y}\mathbf{Y}}, \alpha_w - n + l)} \cdot e^{-\frac{1}{2}(\boldsymbol{\mu}_{\mathbf{Y}} - \boldsymbol{\nu}'_{\mathbf{Y}})^T (N + \alpha_{\mu}) W_{\mathbf{Y}} (\boldsymbol{\mu}_{\mathbf{Y}} - \boldsymbol{\nu}'_{\mathbf{Y}})} e^{-\frac{1}{2} \text{Tr}[R_{\mathbf{Y}\mathbf{Y}} W_{\mathbf{Y}}]}$$

where handily we may simply select the corresponding elements of the posterior parametric matrix R in (A.15).

Finally we marginalise by integrating out $(\mu_{\mathbf{Y}}, W_{\mathbf{Y}})$

$$\begin{aligned}
 p(d^{\mathbf{Y}}) &= \frac{Z(l, N + \alpha_{\mu}, R_{\mathbf{Y}\mathbf{Y}}, N + \alpha_w - n + l)}{(2\pi)^{\frac{lN}{2}} Z(l, \alpha_{\mu}, T_{\mathbf{Y}\mathbf{Y}}, \alpha_w - n + l)} \\
 (A.27) \quad &= \left(\frac{\alpha_{\mu}}{N + \alpha_{\mu}} \right)^{\frac{l}{2}} \frac{\Gamma_l \left(\frac{N + \alpha_w - n + l}{2} \right)}{\pi^{\frac{lN}{2}} \Gamma_l \left(\frac{\alpha_w - n + l}{2} \right)} \frac{|T_{\mathbf{Y}\mathbf{Y}}|^{\frac{\alpha_w - n + l}{2}}}{|R_{\mathbf{Y}\mathbf{Y}}|^{\frac{N + \alpha_w - n + l}{2}}}
 \end{aligned}$$

This result differs from equation (18) of Geiger and Heckerman (2002) in how the matrix elements are selected and differs from equation (24) of Heckerman and Geiger (1995) in how the α_w parameter changes when looking at subsets. Both differences probably derive from the arguments of the Wishart distribution in (A.23) following from Theorem 5.1.4 of Press (1982) itemized above.

B. Simplifying the score. As can be seen in (A.1), the score for each node i involves finding the probability of the data restricted to the node's parents compared to the data restricted to the parents and the node itself. Let \mathbf{P} denote the parent set of size p and \mathbf{Q} the parent set plus the node, of size $p + 1$. Since the difference in size is exactly 1, the ratio of scores involves similar matrix elements and prefactors which can be simplified. Starting with the prefactors:

$$(B.1) \quad \frac{p(d^{\mathbf{Q}})}{p(d^{\mathbf{P}})} = \left(\frac{\alpha_{\mu}}{N + \alpha_{\mu}} \right)^{\frac{1}{2}} \frac{\Gamma \left(\frac{N + \alpha_w - n + p + 1}{2} \right)}{\pi^{\frac{N}{2}} \Gamma \left(\frac{\alpha_w - n + p + 1}{2} \right)} \frac{|T_{\mathbf{Q}\mathbf{Q}}|^{\frac{\alpha_w - n + p + 1}{2}} |R_{\mathbf{P}\mathbf{P}}|^{\frac{N + \alpha_w - n + p}{2}}}{|T_{\mathbf{P}\mathbf{P}}|^{\frac{\alpha_w - n + p}{2}} |R_{\mathbf{Q}\mathbf{Q}}|^{\frac{N + \alpha_w - n + p + 1}{2}}}$$

since the ratios of multivariate gamma functions each leave a single term. The term in front of the ratios of determinants need only be calculated once for each value of $p = 0, \dots, n - 1$ and stored for computational efficiency.

B.1. Simplifying ratios of determinants. The more expensive step is calculating the determinants. For the ratio however one can first partition the larger matrices

$$(B.2) \quad R_{\mathbf{Q}\mathbf{Q}} = \begin{pmatrix} a & \mathbf{b}^T \\ \mathbf{b} & D \end{pmatrix}, \quad D = R_{\mathbf{P}\mathbf{P}}$$

so that

$$(B.3) \quad |R_{\mathbf{Q}\mathbf{Q}}| = a \left| D - \mathbf{b}a^{-1}\mathbf{b}^T \right| = |D| \left(a - \mathbf{b}^T D^{-1} \mathbf{b} \right)$$

and hence

$$(B.4) \quad \frac{|R_{\mathbf{Q}\mathbf{Q}}|^{\frac{N+\alpha_w-n+p+1}{2}}}{|R_{\mathbf{P}\mathbf{P}}|^{\frac{N+\alpha_w-n+p}{2}}} = |D|^{\frac{1}{2}} \left(a - \mathbf{b}^T D^{-1} \mathbf{b} \right)^{\frac{N+\alpha_w-n+p+1}{2}}$$

Due to the difference in powers, $|D|$ still needs to be calculated. Typically, determinants are determined via an LU decomposition, which for the symmetric matrices considered here reduces to the Cholesky decomposition. For a $p \times p$ matrix, this takes $\frac{p^3}{3}$ operations. Although the inverse can be found from the Cholesky decomposition (Krishnamoorthy and Menon, 2011) this is more expensive than taking a second $p \times p$ determinant in the first expression in (B.3). Only for small p would this in turn be notably more efficient than simply calculating $|R_{\mathbf{Q}\mathbf{Q}}|$. In principle one could resort to Coppersmith and Winograd (1990) or faster (Williams, 2012) algorithms for matrix multiplication and inversion to obtain D^{-1} with no overhead on calculating its determinant asymptotically. Here, luckily we can avoid such complications since from the Cholesky decomposition

$$(B.5) \quad D = LL^T, \quad \mathbf{b}^T D^{-1} \mathbf{b} = \mathbf{c}^T \mathbf{c}, \quad \mathbf{c} = L^{-1} \mathbf{b}$$

we know L and can solve $L\mathbf{c} = \mathbf{b}$ for \mathbf{c} using back-substitution with p^2 operations. Using the block partitioning therefore speeds up the calculation of (B.4) by approximately a factor of two.

Using the block partitioning also means we completely avoid matrices when $p = 1$. When $p = 2$ it is much faster to use the explicit form for the determinant and inverse of a 2×2 matrix rather than decompose and back-solve. For $p = 3$ the Laplace expansion may provide the determinant more than twice as quickly as decomposition. In this case there should be a slight additional speed up from using the first expression in (B.3).

These simplifications provide further advantages above the general factor of 2 from the block partitioning, particularly for the sparse DAGs considered in Friedman and Koller (2003) and Grzegorzczak and Husmeier (2008). Since the scoring is likely to be the most expensive part of general MCMC schemes on DAGs (like the structure MCMC of Madigan and York, 1995) the efficiency gained here allows longer chains to be run.

B.2. Diagonal prior parametric matrix. A suggested choice for the prior parametric matrix T , if $\alpha_w > n + 1$, is to approximate the multivariate t -distribution of \mathbf{x} by a normal (Geiger and Heckerman, 2002). Assuming independent normals, leads to

$$(B.6) \quad T = tI_n, \quad t = \frac{\alpha_\mu(\alpha_w - n - 1)}{(\alpha_\mu + 1)}$$

being diagonal (the scaling is to match the covariance of both distributions). Of course, if T is diagonal, then calculating determinants of its subblocks is straightforward. When the elements are equal as here then

$$(B.7) \quad \frac{|T_{\mathbf{Q}\mathbf{Q}}|^{\frac{\alpha_w - n + p + 1}{2}}}{|T_{\mathbf{P}\mathbf{P}}|^{\frac{\alpha_w - n + p}{2}}} = t^{\frac{\alpha_w - n + 2p + 1}{2}}$$

so that this term takes little computational time. Such a choice for T therefore reduces the computation of the score by an additional factor of two.

B.3. *Bias in the score of Heckerman and Geiger (1995).* Finally we compare the terms in (B.1) here to those of Heckerman and Geiger (1995) since that version is likely to have been used in earlier papers employing the BGe score. After simplifying the corresponding ratio of multivariate gamma functions, one obtains

$$(B.8) \quad \frac{p(d^{\mathbf{Q}})}{p(d^{\mathbf{P}})} = \left(\frac{\alpha_\mu}{N + \alpha_\mu} \right)^{\frac{1}{2}} \frac{\Gamma\left(\frac{N + \alpha_w - p}{2}\right)}{\pi^{\frac{N}{2}} \Gamma\left(\frac{\alpha_w - p}{2}\right)} \frac{|T_{\mathbf{Q}\mathbf{Q}}|^{\frac{\alpha_w}{2}} |R_{\mathbf{P}\mathbf{P}}|^{\frac{N + \alpha_w}{2}}}{|T_{\mathbf{P}\mathbf{P}}|^{\frac{\alpha_w}{2}} |R_{\mathbf{Q}\mathbf{Q}}|^{\frac{N + \alpha_w}{2}}}$$

Aside from the difference in the powers of the determinants, the remaining ratio of gamma functions decreases with p unlike the ratio in (B.1). To compare the difference in the prefactor better, we consider their ratio

$$(B.9) \quad g(p) = \frac{\Gamma\left(\frac{N + \alpha_w - n + p + 1}{2}\right)}{\Gamma\left(\frac{\alpha_w - n + p + 1}{2}\right)} \cdot \frac{\Gamma\left(\frac{\alpha_w - p}{2}\right)}{\Gamma\left(\frac{N + \alpha_w - p}{2}\right)}$$

For simplicity, assume that n is odd so that $g = 1$ for $p = \frac{n-1}{2}$. Around this central value

$$(B.10) \quad g(\tilde{p}) = \frac{\Gamma\left(\frac{2N + 2\alpha_w - n + 1}{4} + \frac{\tilde{p}}{2}\right)}{\Gamma\left(\frac{2N + 2\alpha_w - n + 1}{4} - \frac{\tilde{p}}{2}\right)} \cdot \frac{\Gamma\left(\frac{2\alpha_w - n + 1}{4} - \frac{\tilde{p}}{2}\right)}{\Gamma\left(\frac{2\alpha_w - n + 1}{4} + \frac{\tilde{p}}{2}\right)}, \quad \tilde{p} = p - \frac{n-1}{2}$$

Using the recurrence relation for gamma functions

$$(B.11) \quad g(\tilde{p}) = \frac{\left(\frac{2N + 2\alpha_w - n + 1}{4} - \frac{1}{2}\right) \cdots \left(\frac{2N + 2\alpha_w - n + 1}{4} - \frac{\tilde{p}}{2}\right)}{\left(\frac{2\alpha_w - n + 1}{4} - \frac{1}{2}\right) \cdots \left(\frac{2\alpha_w - n + 1}{4} - \frac{\tilde{p}}{2}\right)}, \quad \tilde{p} = 0, \dots, \frac{n+1}{2}$$

we have a simple product of \tilde{p} terms and the inverse for negative \tilde{p} . The exact behaviour of g obviously depends on the dimension n , the parameter

α_w and the number of observations N , but if we assume the latter to be significantly larger than the others ($N \gg n, \alpha_w$) then

$$(B.12) \quad g(\tilde{p}) \sim \left(\frac{N}{2}\right)^{\tilde{p}}, \quad g(p) \sim \left(\frac{N}{2}\right)^{p - \frac{n-1}{2}}$$

with the same behaviour for even n .

Practically this means that the score of Heckerman and Geiger (1995) penalises each node in a DAG with p parents by roughly $\left(\frac{N}{2}\right)^p$ compared to the score in (B.1). Sparse DAGs with few parents are then artificially favoured. Even without assuming N to be much larger than n and α_w , the ratio $g(\tilde{p})$ still increases with \tilde{p} and a bias remains towards sparse DAGs.

References.

- COPPERSMITH, D. and WINOGRAD, S. (1990). Matrix multiplication via arithmetic progressions. *Journal of Symbolic Computation* **9** 251–280.
- FRIEDMAN, N. and KOLLER, D. (2003). Being Bayesian about network structure. A Bayesian approach to structure discovery in Bayesian networks. *Machine Learning* **50** 95–125.
- GEIGER, D. and HECKERMAN, D. (1994). Learning Gaussian networks. In *Tenth Conference on Uncertainty in Artificial Intelligence* 235–243.
- GEIGER, D. and HECKERMAN, D. (2002). Parameter priors for directed acyclic graphical models and the characterization of several probability distributions. *Annals of Statistics* **30** 1412–1440.
- GRZEGORCZYK, M. and HUSMEIER, D. (2008). Improving the structure MCMC sampler for Bayesian networks by introducing a new edge reversal move. *Machine Learning* **71** 265–305.
- HECKERMAN, D. and GEIGER, D. (1995). Learning Bayesian networks: A unification for discrete and Gaussian domains. In *Eleventh Conference on Uncertainty in Artificial Intelligence* 274–284.
- KRISHNAMOORTHY, A. and MENON, D. (2011). Matrix inversion using Cholesky decomposition. Preprint, arXiv:1111.4144.
- MADIGAN, D. and YORK, J. (1995). Bayesian graphical models for discrete data. *International Statistical Review* **63** 215–232.
- PRESS, S. J. (1982). *Applied multivariate analysis*. Krieger publishing company.
- WILLIAMS, V. V. (2012). Multiplying matrices faster than Coppersmith-Winograd. In *Forty-fourth annual ACM symposium on Theory of Computing* 887–898.

JACK KUIPERS
INSTITUT FÜR THEORETISCHE PHYSIK
UNIVERSITÄT REGENSBURG
D-93040 REGENSBURG, GERMANY
E-MAIL: jack.kuipers@ur.de

GIUSI MOFFA
INSTITUT FÜR FUNKTIONELLE GENOMIK
UNIVERSITÄT REGENSBURG
JOSEF ENGERTSTRASSE 9
93053 REGENSBURG, GERMANY
E-MAIL: gusi.moffa@ukr.de

DAVID HECKERMAN
MICROSOFT RESEARCH
1100 GLENDON AVE SUITE PH1
LOS ANGELES CA 90024
E-MAIL: heckerma@microsoft.com